August 2, 2021

# An Introduction to Habana AI Processors

Sree Ganesan

AI Software Product Management @ Habana Labs

www.habana.ai

**habana**
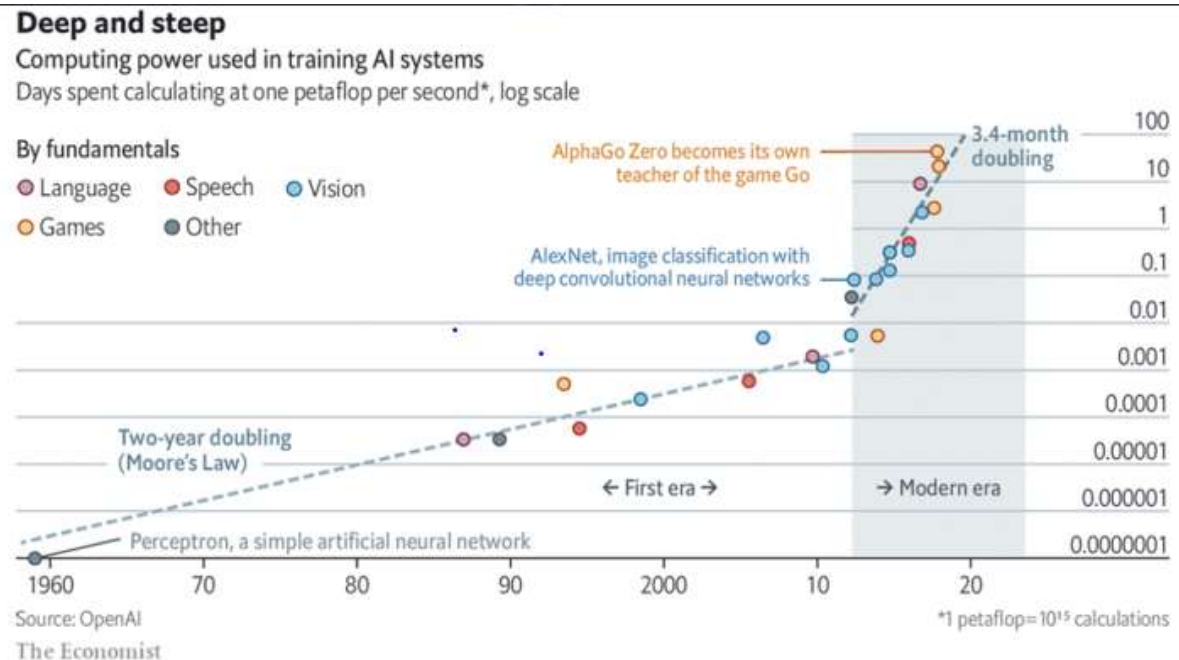An Intel Company

# Agenda

# Demand for compute for ML training doubles every 3.4 months

- **Increasing Complexity**
  - Businesses need higher precision in their model predictions
  - Results in larger and more complex models
  - Requires frequent retraining of models

- **Increasing Costs**
  - Increasing compute power required for frequent training of larger models drives up cost to train
  - Becomes a barrier for innovation and growth



**Deep and steep**
Computing power used in training AI systems
Days spent calculating at one petaflop per second*, log scale

By fundamentals
○ Language  ● Speech  ◎ Vision
○ Games  ● Other

AlphaGo Zero becomes its own teacher of the game Go — 3.4-month doubling

AlexNet, image classification with deep convolutional neural networks

Two-year doubling (Moore's Law)

← First era →   → Modern era

Perceptron, a simple artificial neural network

1960  70  80  90  2000  10  20

Source: OpenAI
The Economist

*1 petaflop=10¹⁵ calculations

**Need for dedicated AI processors to address the compute, memory and communication challenges**

# A little about Habana

- Founded in 2016 to develop purpose-built AI processors

- Launched inference processor in 2018, training processor in 2019

- Acquired by Intel in late-2019

- Fully leveraging Intel's scale, resources and infrastructure

- Accessing Intel ecosystem and customer partnerships

- Delivering aggressive roadmap optimized for AI data center performance and efficiency

**intel**

**&**

**∴habana**®

# Habana's Dedicated Focus: AI Training and Inference

**Training Solution**

**Inference Solution**

GAUDI™

Designed for AI training
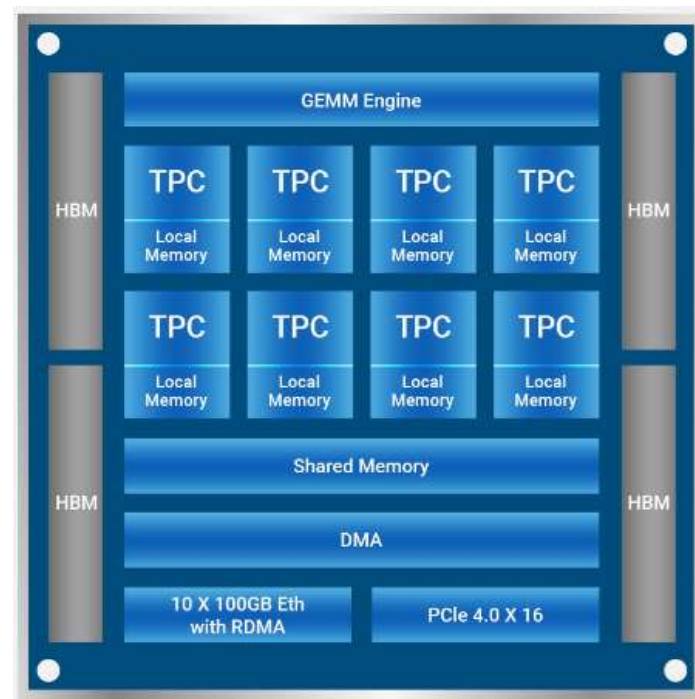efficiency, flexibility and scale

# A new class of AI Training: Habana Gaudi

## Purpose-designed for data center AI Training efficiency

- Cost-efficient AI Training

- Flexibility to ease model migration

- Hardware and software architected for scalability

# Gaudi: Architected for performance and efficiency

• Fully programmable Tensor Processing Cores (TPC) with tools & libraries

• Configurable Matrix Math Engine (GEMM)

• Multi-stage memory hierarchy with 32GB HBM2 memory

• Integrated 10 x 100 Gigabit Ethernet for multi-chip scale-out training

• Delivers higher efficiency than traditional CPUs and GPUs

# Designed for flexible and easy model migration

| Ease of use | Customization | Balanced compute & memory |
|---|---|---|
| Integrated with TensorFlow and PyTorch; minimal code changes to get started<br>➜ SynapseAI maps model topology onto Gaudi devices | SynapseAI TPC SDK facilitates development of custom kernels | 32GB HBM2 memories similar to GPUs, so existing DL models will fit into Gaudi memory |
| Developers can enjoy the **same abstraction** they are accustomed to today | Developers can **customize** models to extract best performance | Developers can spend **less effort** to port their models to Gaudi |

# Designed for Scaling Efficiency

The industry's FIRST:
Native integration of 10 x 100 Gigabit Ethernet RoCE ports onto every Gaudi

- Eliminates network bottlenecks

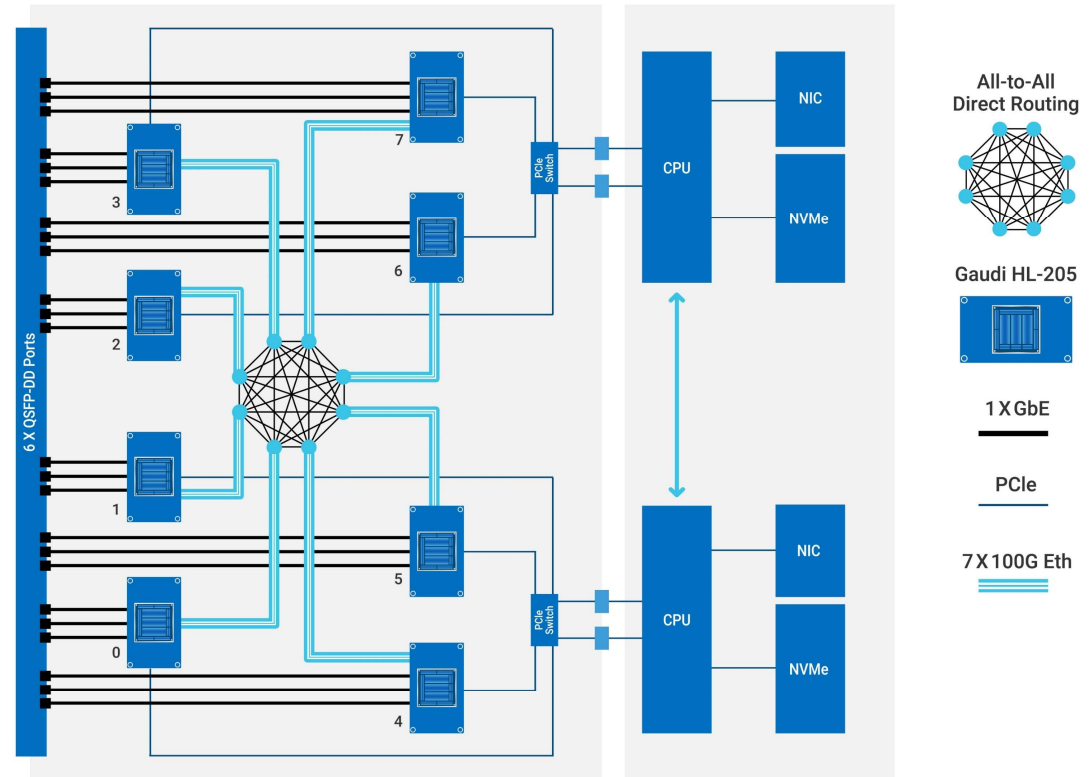- Standard Ethernet inside the server and across nodes

- Eliminates lock-in with proprietary interfaces

- Lowers total system cost and power by reducing discrete components

# Scaling within a Gaudi Server

- 8 Gaudi OCP OAM cards

- 24 x 100GbE RDMA RoCE for scale-out

- Non-blocking, all-2-all internal interconnect across Gaudi AI processors

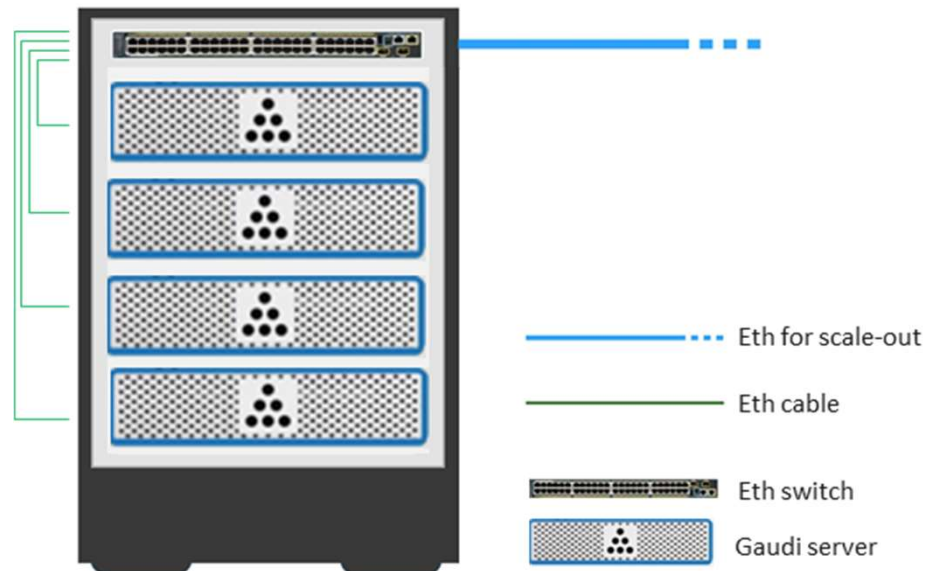- Separate PCIe ports for external Host CPU traffic



Example of Integrated Server with eight Gaudi AI processors, two Xeon CPU and multiple Ethernet Interfaces

# Rack and Pod Level Scaling

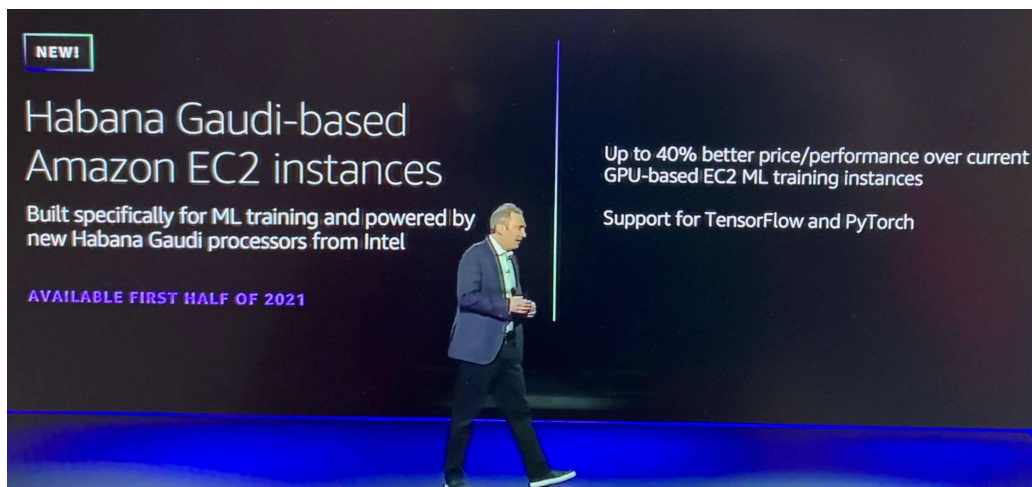Easily build rack and pod-scale training systems with off-the-shelf standard ethernet switches

Eth for scale-out

Eth cable

Eth switch

Gaudi server

Example of rack configuration with four Gaudi servers (eight Gaudi processors per server) connected to a single Ethernet switch

# Gaudi-based
# Amazon EC2 AI Training Instances

# Gaudi-Based AWS EC2 Instances Coming Soon



*"The new EC2 instances will leverage up to 8 Gaudi accelerators and deliver **up to 40% better price/performance** than current GPU-based EC2 instances for training DL models."*

*Andy Jassy, re:Invent 2020*

- Amazon's first non-GPU instances based on Habana Gaudi AI processors
- Improved cost-efficiency makes AI Training accessible to more customers

# Gaudi-Based EC2 Instances



Benefit from full stack of Amazon EC2 services:

- AWS DLAMI, DLC for Gaudi

- AWS ECS and EKS orchestration for containerized applications

- Integration with Amazon SageMaker

- Efficient scaling across multiple Gaudi-based EC2 Instances

# On-Premise Solution

# Partnering with Supermicro

## Solutions available now for on-premises customers

## Featured Servers:

- Supermicro X12 Gaudi® AI Training System
  - Eight Gaudi HL-205 AI processors
  - Dual-socket 3rd Gen Intel® Xeon® Scalable processors
  - https://www.supermicro.com/en/products/system/AI/4U/SYS-420GH-TNGR

- Supermicro SuperServer 4029GP-T
  - Eight Goya™ HL-100 PCIe cards
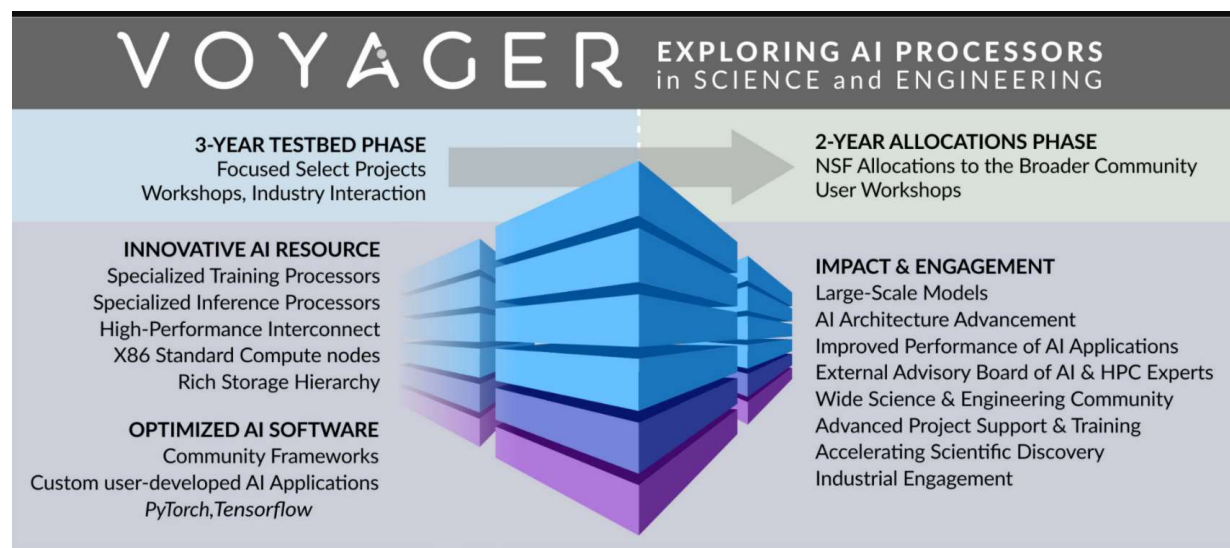  - Dual-socket 2nd Gen Intel® Xeon® Scalable processors
  - https://www.supermicro.com/en/products/system/4U/4029/SYS-4029GP-TRT.cfm

Supermicro
X12 Gaudi
AI Training System

# Habana AI to power SDSC's Voyager Research Program

## 336 Gaudi Training accelerators with native RoCE scaling and 16 Goya Inference processors

- Voyager to go into service Fall 2021
- Funded by $5M grant from National Science Foundation
  - Matching funds targeting community support and operation
- AI research conducted across range of science and engineering domains
  - Astronomy, climate sciences, chemistry, particle physics,
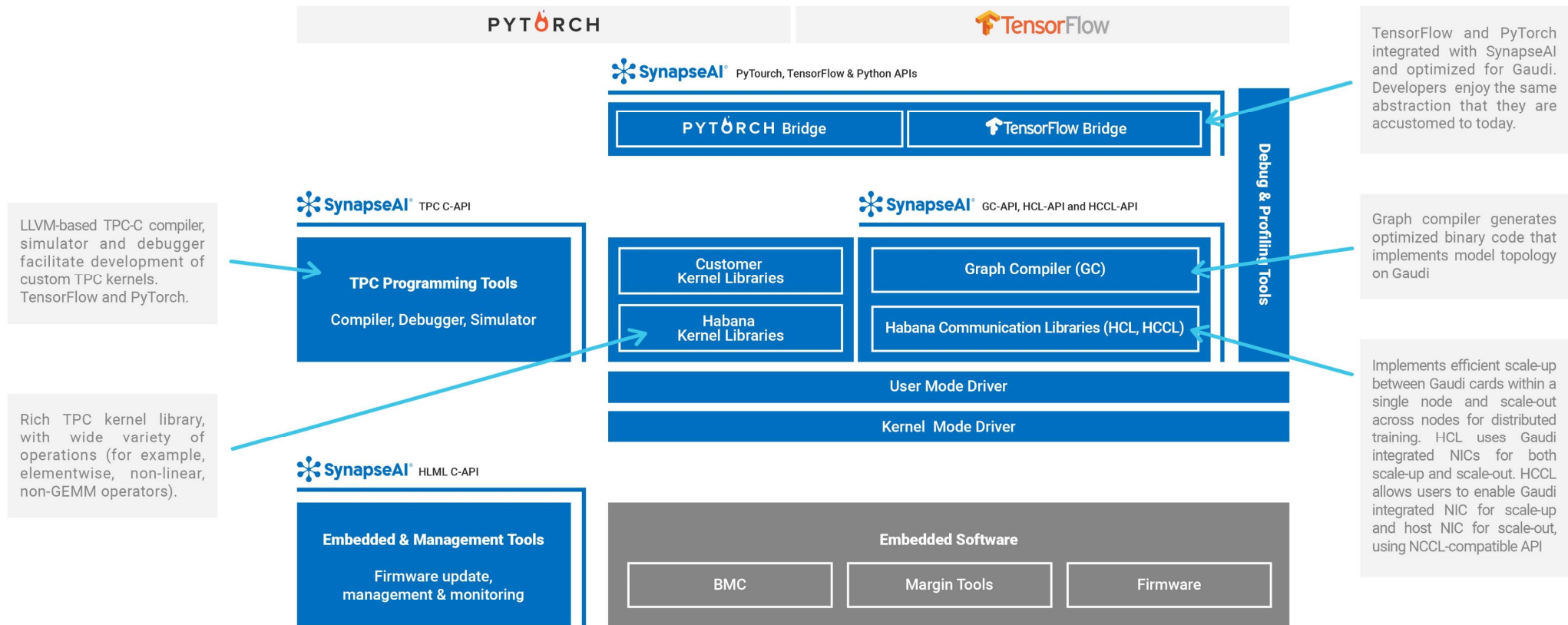- Announced by SDSC in July 2020, more information here.



VOYAGER — EXPLORING AI PROCESSORS in SCIENCE and ENGINEERING

**3-YEAR TESTBED PHASE**
Focused Select Projects
Workshops, Industry Interaction

**2-YEAR ALLOCATIONS PHASE**
NSF Allocations to the Broader Community
User Workshops

**INNOVATIVE AI RESOURCE**
Specialized Training Processors
Specialized Inference Processors
High-Performance Interconnect
X86 Standard Compute nodes
Rich Storage Hierarchy

**OPTIMIZED AI SOFTWARE**
Community Frameworks
Custom user-developed AI Applications
PyTorch, Tensorflow

**IMPACT & ENGAGEMENT**
Large-Scale Models
AI Architecture Advancement
Improved Performance of AI Applications
External Advisory Board of AI & HPC Experts
Wide Science & Engineering Community
Advanced Project Support & Training
Accelerating Scientific Discovery
Industrial Engagement

# Software,
# Resources & Support

# SynapseAI Software Suite

- Train deep learning models on Gaudi with minimal code changes

- Natively integrated with TensorFlow & PyTorch

- Reference models, kernel libraries, software and docs available on GitHub

- Advanced users can write their own custom software kernels

# Software Suite Detail



PYT🔥RCH

TensorFlow

SynapseAI® PyTourch, TensorFlow & Python APIs

PYT🔥RCH Bridge

TensorFlow Bridge

TensorFlow and PyTorch integrated with SynapseAI and optimized for Gaudi. Developers enjoy the same abstraction that they are accustomed to today.

SynapseAI® TPC C-API

SynapseAI® GC-API, HCL-API and HCCL-API

LLVM-based TPC-C compiler, simulator and debugger facilitate development of custom TPC kernels. TensorFlow and PyTorch.

**TPC Programming Tools**

Compiler, Debugger, Simulator

Customer Kernel Libraries

Graph Compiler (GC)

Habana Kernel Libraries

Habana Communication Libraries (HCL, HCCL)

Debug & Profiling Tools

Graph compiler generates optimized binary code that implements model topology on Gaudi

User Mode Driver

Kernel Mode Driver

Rich TPC kernel library, with wide variety of operations (for example, elementwise, non-linear, non-GEMM operators).

SynapseAI® HLML C-API

**Embedded & Management Tools**

Firmware update, management & monitoring

**Embedded Software**

BMC

Margin Tools

Firmware

Implements efficient scale-up between Gaudi cards within a single node and scale-out across nodes for distributed training. HCL uses Gaudi integrated NICs for both scale-up and scale-out. HCCL allows users to enable Gaudi integrated NIC for scale-up and host NIC for scale-out, using NCCL-compatible API

# TensorFlow integration with SynapseAI



SynapseAI receives a computational graph of the model from the framework

It identifies subgraphs (blue nodes) that can be accelerated by Gaudi

The rest of the graph runs on CPU (yellow node)

The original graph is modified to replace the Gaudi subgraphs with encapsulated nodes (blue)

The framework runtime executes the modified graph

For each encapsulated node, SynapseAI generates optimized binary code that runs on Gaudi

# Getting Started with TensorFlow on Gaudi

```
import tensorflow as tf

from TensorFlow.common.library_loader import load_habana_module
load_habana_module()

(x_train, y_train), (x_test, y_test) = tf.keras.datasets.mnist.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0

model = tf.keras.models.Sequential([
                    tf.keras.layers.Flatten(input_shape=(28, 28)),
                    tf.keras.layers.Dense(10),
])
loss = tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True)
optimizer = tf.keras.optimizers.SGD(learning_rate=0.01)

model.compile(optimizer=optimizer, loss=loss, metrics=['accuracy'])

model.fit(x_train, y_train, epochs=5, batch_size=128)
model.evaluate(x_test, y_test)
```

Load the Habana libraries needed to use Gaudi aka **HPU** device.

Once loaded, the **HPU device** is registered in TensorFlow and prioritized over CPU.

When an Op is available for both CPU and HPU, the Op is assigned to the HPU.

When an Op is not supported on HPU, it runs on the CPU

# Habana Developer Platform

# Habana Developer Resources

# Habana's Developer Documentation



https://docs.habana.ai

# Habana Developer Software Vault



## https://vault.habana.ai

### Update Your Software

Download the latest SynapseAI(R) Software, including Habana's graph compiler and runtime, TPC kernel library, firmware and drivers, and tools. These components are needed to update an existing system to the latest drivers and Firmware. For more information on how to install this content, please refer to the Installation Guide.

| Name | Description | Download |
|---|---|---|
| Habanalabs-graph | Installs the Graph Compiler and the run-time. | Ubuntu18.04 Ubuntu20.04 AmazonLinux2 Centos7.5 |
| habanalabs-thunk | Installs the thunk library. | Ubuntu18.04 Ubuntu20.04 AmazonLinux2 Centos7.5 |
| habanalabs-dkms | Installs the PCIe driver. | Ubuntu18.04 Ubuntu20.04 AmazonLinux2 Centos7.5 |
| habanalabs-fw-tools | Installs various Firmware embedded tools (hlml, hl-smi, etc). | Ubuntu18.04 Ubuntu20.04 AmazonLinux2 Centos7.5 |
| habanalabs-aeon | Installs synapse level demo's data loader. | Ubuntu18.04 Ubuntu20.04 AmazonLinux2 Centos7.5 |
| habanalabs-qual | Installs the qualification application package. | Ubuntu18.04 Ubuntu20.04 AmazonLinux2 |

# Habana GitHub

# Software Installation and Deployment

Setup_and_Install repository on Habana GitHub provides instructions on how to setup your environment with the SynapseAI software stack

**SynapseAI Orchestration**
(Kubernetes Gaudi plugin, Kubeflow mpi-operator)

**SynapseAI TensorFlow Container Image**
(TensorFlow frontend, horovod, open-mpi)

**SynapseAI Base Installation Image**
(OS, Gaudi linux kernel driver, user mode driver, graph compiler, HCL/HCCL & embedded tools)

**Gaudi Server**

Gaudi-optimized Docker container images with all necessary dependences*

**Official releases publicly available on Habana Vault**

| | |
|---|---|
| Orchestration | Kubernetes (1.19) |
| Frameworks | TensorFlow2 and PyTorch |
| Operating Systems | Ubuntu 18.04 and 20.04 |
| Container Runtimes | Docker (Docker CE version 18.09) |
| Distributed Training Schemes | TensorFlow with Horovod and tf.distribute PyTorch distributed (native) |

*Habana GitHub will have repository with Dockerfiles to "build your own" Docker images*

# Gaudi Reference Models – July 2021



- Gaudi reference model roadmap on Habana GitHub

- Scripts and detailed instructions to enable the reference models on Gaudi available on the Model-References repository

# Habana's Developer Forum

https://forum.habana.ai

There are many exciting opportunities for deep learning in scientific research.

Habana Labs invites you to explore the possibilities with our Gaudi AI Training Processor!

March 2021

# THANK YOU

www.habana.ai